

日 本 国 特 許 庁
JAPAN PATENT OFFICE

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日 2 0 0 3 年 4 月 1 日
Date of Application:

出 願 番 号 特 願 2 0 0 3 - 0 9 7 5 7 3
Application Number:
[ST. 10/C] : [J P 2 0 0 3 - 0 9 7 5 7 3]

願 人 株式会社日立製作所
Applicant(s):

2 0 0 3 年 9 月 1 8 日

特許庁長官
Commissioner,
Japan Patent Office

今 井 康 夫

【書類名】 特許願

【整理番号】 K03002011A

【あて先】 特許庁長官殿

【国際特許分類】 G06F 9/45

【発明者】

 【住所又は居所】 神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所 システム開発研究所内

 【氏名】 本川 敬子

【発明者】

 【住所又は居所】 神奈川県川崎市麻生区王禅寺 1 0 9 9 番地 株式会社日立製作所 システム開発研究所内

 【氏名】 久島 伊知郎

【発明者】

 【住所又は居所】 神奈川県横浜市戸塚区戸塚町 5 0 3 0 番地 株式会社日立製作所 ソフトウェア事業部内

 【氏名】 伊藤 信一

【特許出願人】

 【識別番号】 000005108

 【氏名又は名称】 株式会社日立製作所

【代理人】

 【識別番号】 100075096

 【弁理士】

 【氏名又は名称】 作田 康夫

【手数料の表示】

 【予納台帳番号】 013088

 【納付金額】 21,000円

【提出物件の目録】

 【物件名】 明細書 1

 【物件名】 図面 1

【物件名】 要約書 1
【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 階層メモリ向け最適化処理を備えたコンパイラおよびコード生成方法

【特許請求の範囲】

【請求項 1】

計算機システムと協働することで、複数のメモリ階層を備えたアーキテクチャ上で実行するためのオブジェクトプログラムを生成するコンパイラであって、

対象プログラムが実行時に主としてどのメモリ階層上のデータを参照するかを指示するオプションまたは指示文を解釈するステップと、

上記指示されたメモリ階層向けの最適化処理を行うステップとを有することを特徴とするコンパイラ。

【請求項 2】

請求項 1 のコンパイラであって、

前記メモリ階層向けの最適化処理として、メモリをアクセスする命令に対して指定されたメモリ階層に応じたメモリレイテンシを求め、求めたレイテンシに応じた最適化処理を行うことを特徴とするコンパイラ。

【請求項 3】

請求項 1 のコンパイラであって、

メモリ階層向けの最適化処理として、メモリをアクセスする命令に対して指定されたメモリ階層に応じて、ループ交換またはループ展開またはループタイリングのループ変換方法を決定することを特徴とするコンパイラ。

【請求項 4】

計算機システムおよび該計算機システムと協働して複数のメモリ階層を備えたアーキテクチャ上で実行するためのオブジェクトプログラムを生成するコンパイラにより遂行されるオブジェクトプログラム生成方法であって、

対象プログラムが実行時に主としてどのメモリ階層上のデータを参照するかを指示するオプションまたは指示文を解釈するステップと、

上記指示されたメモリ階層向けの最適化処理を行うステップとを有することを特徴とするコード生成方法。

【請求項 5】

請求項 4 のコード生成方法であって、

前記メモリ階層向けの最適化処理として、メモリをアクセスする命令に対して指定されたメモリ階層に応じたメモリレイテンシを求め、求めたレイテンシに応じた最適化処理を行うことを特徴とするコード生成方法。

【請求項 6】

請求項 4 のコード生成方法であって、

メモリ階層向けの最適化処理として、メモリをアクセスする命令に対して指定されたメモリ階層に応じて、ループ交換またはループ展開またはループタイリングのループ変換方法を決定することを特徴とするコード生成方法。

【請求項 7】

複数のメモリ階層を備えたアーキテクチャ上で実行するプログラムに対して、対象プログラムが実行時に主としてどのメモリ階層上のデータを参照するかを指示するオプションまたは指示文を解釈することを特徴とするコンパイラ。

【請求項 8】

請求項 1 のコンパイラを格納した記憶媒体。

【発明の詳細な説明】**【 0 0 0 1 】****【発明の属する技術分野】**

本発明は、計算機の利用技術において、オブジェクトプログラムの実行時間を削減するコンパイル方法に関する。特に、複数のメモリ階層を備えたアーキテクチャ向けの最適化指示方法に関する。

【 0 0 0 2 】**【従来の技術】**

マイクロプロセッサの速度向上に伴って主記憶アクセスのレイテンシが増大している。多くのプロセッサでは、主記憶よりもアクセスが速く比較的小容量のキャッシュ記憶を備えており、さらにキャッシュ記憶を、1 次キャッシュ、2 次キャッシュというように階層化することにより、なるべくレイテンシの小さいメモリ階層のデータをアクセスすることで、レイテンシの大きなメモリ階層のアクセ

ス回数を削減している。すなわちメモリアクセス命令は、1次キャッシュヒット時には短いレイテンシで1次キャッシュをアクセスし、1次キャッシュミス時には、2次キャッシュヒットならば2次キャッシュをアクセスし、全てのキャッシュ階層のミス時にのみ主記憶をアクセスする。

【0003】

キャッシュミス時のレイテンシを隠す方法としては、ロード命令と、ロードしたデータを使用する命令の距離を十分離すような命令スケジューリングがあげられる。命令間を何サイクル離すかを定める基準としては、メモリレイテンシを用いる。

【0004】

キャッシュミス時の主記憶参照のレイテンシを隠す別の手法として、例えば非特許文献1などに述べられているように、主記憶からキャッシュへデータを先行的に移動するプリフェッチ命令をプロセッサに用意し、コンパイラによってプログラム中にプリフェッチ命令を挿入する、ソフトウェア・プリフェッチ（プリフェッチ最適化）と呼ばれる方法がある。プリフェッチ命令を利用すれば、後続ループ繰り返しで参照するデータを予め主記憶からキャッシュへ移動しておきながら、同時に別の演算を行うことができ、主記憶参照のレイテンシを隠すことができる。

【0005】

この方法では、ループ内のデータ参照に対してプリフェッチ命令を生成する場合、そのループの1回の繰り返しにかかる実行サイクル数Cを見積もり、データをメモリからキャッシュに移動するのに要するサイクル数（メモリレイテンシ）LをCで割った値 $\alpha = \text{CEIL}(L/C)$ （CEILは小数点以下の切り上げを表す記号とする）を計算し、 α 回後の繰り返しで参照するデータを予めプリフェッチしておくことにより、Lサイクル後にそのデータを参照するときにはデータが既にキャッシュに到着しているのでキャッシュヒットとなり、プログラム実行が高速化される。1次キャッシュへのプリフェッチを行う場合、データが既に1次キャッシュ上にあればプリフェッチは不要である。また、データが2次キャッシュにデータがある場合にはLは2次キャッシュから1次キャッシュへ移動するのに要するサ

イクル数、データが主記憶上にある場合には主記憶から 1 次キャッシュへ移動するのに要するサイクル数を使用した方がよい。しかし、通常、データがどのメモリ階層にあるかは不明なため、主記憶上にあると仮定した処理を行う。

【0 0 0 6】

データ局所性を高めるプログラム変換により、キャッシュミスの回数を削減する最適化も知られている。具体的な変換としては、ループタイリング、ループ交換、ループ展開があげられる。

【0 0 0 7】

ループタイリングとは、多重ループ内で参照されるデータがリユースを持つ場合に、一旦キャッシュ上にロードされたデータが他のデータの参照によってキャッシュから追い出される前に再度参照されるようにすることを目的としたループ変換である。ループタイリングについては、非特許文献 2 に記載されている。

【0 0 0 8】

メモリ参照パターンの最適化を目的としたループ交換やループ展開については、非特許文献 3 に記載されている。

【0 0 0 9】

以上で述べたようなレイテンシ隠蔽最適化やデータ局所化は、メモリ参照にかかるサイクル数や、キャッシュサイズといった、対象とするターゲットマシンの属性に依存する情報を必要とする。コンパイラでは通常、対象とするマシンの情報を内部情報として保持しているが、これらの情報をユーザが指示する方法もある。非特許文献 4 では'-mslatency=N' (Nは正定数) というオプション指定により、メモリからの読み出しにかかるサイクル数がNであることを指定できることが記載されている。非特許文献 5 では'-qcache' オプションにより階層キャッシュのレベル毎にキャッシュサイズやラインサイズ、連想数が指定できることが記載されている。

【0 0 1 0】

【非特許文献 1】 Todd C. Mowry 他：Design and Evaluation of Compiler Algorithm for Prefetching, Architectural Support for Programming Languages and Operating Systems, pp.62-73, 1992年

【非特許文献 2】 Michael Edward Wolf, 'Improving Locality and Parallelism in Nested Loops', Technical Report: CSL-TR-92-538, 1992年

【非特許文献 3】 Kevin Dowd, 'High Performance Computing', O'Reilly & Associates, Inc., 11.1節

【非特許文献 4】 日立, '(HI-UX/MPP for SR8000) Optimizing FORTRAN90 User's Guide', 6.2節

【非特許文献 5】 IBM, 'XL Fortran for AIX User's Guide Version 7 Release 1', 第5章

【 0 0 1 1 】

【発明が解決しようとする課題】

従来の技術の項で述べたレイテンシ隠蔽最適化やデータ局所化は、プログラム実行時にデータがどのメモリ階層に存在するかにより最適化の方法は異なる。

【 0 0 1 2 】

例えば、命令スケジューリングにおいては、ロード命令と使用命令の間の距離が大きくなれば、使用レジスタ数が増大する。またプリフェッチ最適化においては、プリフェッチ命令のタイミングが早過ぎると、使用命令が実行される前に再びキャッシュから追い出されてしまう場合がある。したがって、これらの最適化で仮定するメモリレイテンシが大き過ぎると最適化の効果が十分に得られない。すなわちL2キャッシュヒットのデータに対して主記憶レイテンシを用いた最適化を適用すると、L2キャッシュレイテンシを用いた場合に比べて実行性能が低下する可能性がある。しかし、従来の技術ではロード命令の対象データがどのメモリ階層に存在するかが不明であるため、主記憶レイテンシを仮定せざるを得ない問題がある。

【 0 0 1 3 】

また、データ局所性最適化では、複雑なループ構造に変換したためにループ実行のオーバーヘッドが増加し実行性能が低下する可能性がある。ループ内で参照されるデータが主としてキャッシュミスを起こす場合にはキャッシュミス回数の削減により適用効果が得られるが、キャッシュヒット時にはキャッシュミス回数の削減効果はないため、適用しない方がよい。従来はデータがキャッシュヒット

するかどうかがわからないため、キャッシュヒット時にもループ変換を適用していた。このため実行性能が低下する場合があるという問題があった。

【 0 0 1 4 】

【課題を解決するための手段】

上記の課題を解決するために、本発明では、計算機システムと協働することで、複数のメモリ階層を備えたアーキテクチャ上で実行するためのオブジェクトプログラムを生成するコンパイラに、対象プログラムが実行時に主としてどのメモリ階層上のデータを参照するかを指示するオプションまたは指示文を解釈するステップと、上記指示されたメモリ階層向けの最適化処理を行うステップとを持たせた。

メモリ階層向けの最適化処理としては、メモリをアクセスする命令に対して指定されたメモリ階層に応じたメモリレイテンシを求め、求めたレイテンシに応じた最適化処理を行うステップ／またはメモリをアクセスする命令に対して指定されたメモリ階層に応じて、ループ交換またはループ展開またはループタイリングのループ変換方法を決定するステップを持たせた。

【 0 0 1 5 】

図 1 は本発明の概要を示す。図 1 においてソースプログラム TEST1.f の実行時に参照されるデータが主として L2 キャッシュデータであり、TEST2.f は主記憶データだとすると、これら 2 つのプログラムに対する最適化方法は異なるべきである。

【 0 0 1 6 】

しかし、従来はどちらも主記憶データを仮定した最適化を適用していた。プログラム内で参照されるデータが主としてアクセスするメモリ階層は、データサイズやループサイズなどによって決まる。これらはコンパイラの静的な解析では不明であることが多いためである。

【 0 0 1 7 】

本発明では、プログラム実行時にデータが主としてどのメモリ階層に属するかを指定する手段を設けることにより、コンパイラがメモリ階層の指定を解析し（1 0 1）、指定されたメモリ階層に応じた最適化を実施する（1 0 3、1 0 4、

1 0 5) ことで、より進んだ最適化が実施されたオブジェクトプログラムが生成できる。メモリ最適化では、指定されたメモリ階層のレイテンシやサイズなどの属性を用いてコード変換を実施する。図 1 の例ではTEST1.fに対してはL2データ向け最適化(1 0 4)を適用し、TEST2.fに対しては主記憶データ向け最適化(1 0 5)を適用する。

【0 0 1 8】

【発明の実施の形態】

以下、図面を用いて本発明の実施の形態について説明する。

【0 0 1 9】

図 2 は、本発明によるコンパイラが稼動する計算機システムの構成図である。この計算機システムは、CPU 2 0 1、ディスプレイ装置 2 0 2、キーボード 2 0 3、主記憶装置 2 0 4、および外部記憶装置 2 0 5 より構成されている。キーボード 2 0 3 により、ユーザからのコンパイラ起動命令を受け付ける。コンパイラ終了メッセージやエラーメッセージは、ディスプレイ装置 2 0 2 に表示される。外部記憶装置 2 0 5 には、ソースプログラム 2 0 6 とオブジェクトプログラム 2 0 7 が格納される。主記憶装置 2 0 4 には、コンパイラ 2 0 8、コンパイル過程で必要となる中間コード 2 0 9 とループ表 2 1 0 が格納される。コンパイル処理はCPU 2 0 1 がコンパイラプログラム 2 0 8 を実行することにより行われる。なお、CPU 2 0 1 は内部にフェッチ／デコードユニット、実行ユニット等を含む処理ユニット 2 0 1 1 からのアクセスレイテンシの異なるレベル 1 (L 1) キャッシュ 2 0 1 2 とレベル 2 (L 2) キャッシュ 2 0 1 3 を持ち、処理ユニットによる処理を高速に実行するためのメモリ階層を構成している。また、以降の説明では、コンパイラに記述されたコードを解釈して実行するCPU 2 0 1 による処理をコンパイラの処理／処理手順と言うものとする。

【0 0 2 0】

図 3 に、コンパイル時にメモリ階層指示を行うために、計算機システムが入力 J を受け付けるコンパイラ起動命令の例を示す。'f90' はFortranコンパイラの起動コマンドの例、test.f はソースプログラムの指定の例、'-O' はコンパイラオプションの一例である。本実施例ではメモリ階層指示オプションは'-data' で表し

、プログラム中のデータが主としてどのメモリ階層に存在するかを指定する。メモリ階層としては、L1キャッシュを表す'L1'、L2キャッシュを表す'L2'、主記憶を表す'MEM'を指定できる。'-data=L2' (3 0 1) では、データが主としてL2キャッシュに存在することの指示をしている。

【0 0 2 1】

図3の例ではコンパイラ起動命令中でメモリ階層指示を行う例を示すが、メモリ階層指示は、ソースプログラム中の指示文として記述することもできる。図4にメモリ階層指示が付加されたソースプログラムの例を示す。本例では、ループ単位に指示文を指定できるものとする。指示文4 0 1は'D0 10'のループ内のデータが主としてL1キャッシュヒットとなることの指示をしている。指示文4 0 2は'D0 20'のループ内のデータが主として主記憶をアクセスすることの指示をしている。指示文4 0 3は'D0 30'のループ内のデータが主としてL2キャッシュヒット (L1キャッシュミス) となることの指示をしている。

【0 0 2 2】

図5に、図2のシステムで稼動するコンパイラ2 0 8の処理手順を示す。コンパイラの処理は、構文解析5 0 1、メモリ最適化5 0 2、レジスタ割り付け5 0 3、コード生成5 0 4の順で行う。

【0 0 2 3】

構文解析5 0 1では、ソースプログラム2 0 6を入力として構文解析、およびループ解析を行い、中間コード2 0 9とループ表2 1 0を出力する。また本ステップにおいて、本発明の特徴であるメモリ階層指示文やメモリ階層指示オプションを解析し、ループ表に登録する。ループ表2 1 0については後で説明する。メモリ最適化5 0 2は、本発明の特徴となるステップであるので、後で詳細に説明する。レジスタ割り付け5 0 3では中間コード2 0 9の各ノードへのレジスタ割り付けを行う。コード生成5 0 4では、中間コード2 0 9をオブジェクトプログラム2 0 7に変換し、出力する。

【0 0 2 4】

図6にループ表2 1 0の例を示す。図には、ループを識別するためのループ番号6 0 1のほかは、本発明の特徴であるデータ指示フィールド6 0 2だけを示し

ている。ループ表は構文解析処理 5 0 1 の処理で生成されるが、その生成処理であるデータ指示フィールド設定の処理手順を図 1 3 に示す。

【 0 0 2 5 】

本処理は、ループを順に辿り処理する。ステップ 1 3 0 1 では、未処理のループがあるかどうかを判定し、あればステップ 1 3 0 2 で未処理ループを 1 つ取り出す。ステップ 1 3 0 3 ではコンパイル対象のソースプログラム中に図 4 に示したようなループへのメモリ階層指示である 'data' 指示文があるかどうかを調べる。あればステップ 1 3 0 7 で 'data' 指示文のオペランドをループ表に設定する。図 4 の例では D010 のループに対しては 'LI' のように設定する。data 指示文がないときにはステップ 1 3 0 4 へ進み、コンパイラ起動命令にメモリ階層指示オプションである '-data' オプションの指定があるかどうかを調べる。オプション指定がある場合はステップ 1 3 0 5 へ進み、オプションの値を設定する。オプション指定がなければステップ 1 3 0 6 へ進み、「指定なし」を設定する。

【 0 0 2 6 】

以下、メモリ最適化 5 0 2 の詳細な例を示す。

【 0 0 2 7 】

最適化処理の第 1 の例はメモリ最適化 5 0 2 において命令スケジューリングを行うものとする。命令スケジューリングの処理手順を図 7 に示す。

【 0 0 2 8 】

ステップ 7 0 1 で、DAG(directed acyclic graph)を作成する。図 8 に $A(I)*B(I)+C(I)$ に対応した DAG の例を示す。DAG の各ノードは中間コード上の命令に対応している。ノード間のエッジは、実行順序の制約を表す。例えばノード 8 0 3 の mul 演算はノード 8 0 1 およびノード 8 0 2 の load の結果をオペランドとするので、ノード 8 0 1 とノード 8 0 2 は、ノード 8 0 3 よりも先に実行しなければならない。この関係をノード 8 0 1 からノード 8 0 3 へのエッジ、およびノード 8 0 2 からノード 8 0 3 へのエッジで表す。

【 0 0 2 9 】

ステップ 7 0 2 で、ステップ 7 0 1 で作成した DAG のエッジ上にレイテンシを設定する。本ステップの処理手順の詳細を図 9 に示す。図 9 の各ステップについ

て説明する。

【0 0 3 0】

本処理は、DAG上のエッジを辿り、順に処理する。ステップ9 0 1では未処理のエッジがあるかどうかを判定し、なければ処理を終了する。あればステップ9 0 2へ進み、未処理のエッジを取り出して、処理対象とする。ステップ9 0 3では、エッジの始点ノードがロード命令かどうかを調べる。ロード命令でなければステップ9 1 0へ進み、始点ノードの演算に対応するレイテンシの値をエッジ上に設定して、このエッジの処理を終了し、ステップ9 0 1へ戻る。

【0 0 3 1】

ステップ9 0 4では、ノードの所属するループを調べ、ループ表によりdata指示を調べる。ステップ9 0 5では、data指示がL1かどうかを調べる。そうであればステップ9 0 9へ進み、L1キャッシュのレイテンシの値をエッジ上に設定してステップ9 0 1へ戻る。data指定がL1でなければステップ9 0 6へ進み、data指示がL2かどうかを調べる。そうであればステップ9 0 8へ進み、L2キャッシュのレイテンシの値をエッジ上に設定してステップ9 0 1へ戻る。そうでない場合、即ちdata指示がMEMまたは指定なしの場合にはステップ9 0 7へ進み、主記憶アクセスのレイテンシを設定してステップ9 0 1へ戻る。

【0 0 3 2】

図8は各ロードノードに対応するdata指示がL2で、L2レイテンシが10サイクルの場合の例を示している。

【0 0 3 3】

本実施例では、ロード命令の対象データがL1キャッシュ、L2キャッシュ、主記憶のいずれに存在しているかに応じてDAGエッジのレイテンシを設定することにより、ロード命令に対してより正確なレイテンシを仮定した命令スケジューリングを適用できるという効果がある。

【0 0 3 4】

次に、本発明によるメモリ最適化5 0 2において、プリフェッチ最適化を行う例を説明する。プリフェッチ最適化の処理手順を図1 0に示す。

【0 0 3 5】

ステップ 1 0 0 1 では、未処理のループがあるかどうかを判定し、あればステップ 1 0 0 2 で未処理ループを 1 つ取り出す。ステップ 1 0 0 3 では取り出したループの data 指示をループ表 2 1 0 により調べ、data 指示が L1 ならば、ステップ 1 0 0 1 へ戻る。これは、このループのデータが既に L1 キャッシュ上にあるため、プリフェッチ最適化の効果がないと判断したためである。ステップ 1 0 0 4 では data 指示が L2 かどうかを調べる。L2 ならばステップ 1 0 0 6 へ進み、L に L2 キャッシュのレイテンシを設定する。そうでない場合、即ち data 指示が MEM または指定なしの場合は、ステップ 1 0 0 5 へ進み、L に主記憶レイテンシを設定する。

【 0 0 3 6 】

ステップ 1 0 0 7 ではループ繰り返しあたりの実行サイクル数を求め、C に設定する。ステップ 1 0 0 8 では (L/C) を整数に切り上げた値を N に設定する。ステップ 1 0 0 9 では N イタレーション後のプリフェッチコードを生成し、ループ内に挿入する。

【 0 0 3 7 】

本実施例では L2 キャッシュにあるかどうかに応じて主記憶レイテンシを設定することにより、より適切なプリフェッチ距離によるプリフェッチコードを生成する効果がある。また、L1 キャッシュにデータがある場合には不要なプリフェッチコードを生成しない効果がある。

【 0 0 3 8 】

次に、本発明によるメモリ最適化 5 0 2 において、タイリングを行う例について説明する。タイリングの処理手順を図 1 1 に示す。

【 0 0 3 9 】

ステップ 1 1 0 1 では、未処理のループがあるかどうかを判定し、あればステップ 1 1 0 2 で未処理ループを 1 つ取り出す。ステップ 1 1 0 3 では取り出したループの data 指示をループ表 2 1 0 により調べ、data 指示が L1 ならば、ステップ 1 1 0 1 へ戻る。これは、このループのデータが既に L1 キャッシュ上にあるため、タイリングの効果がないと判断したためである。ステップ 1 1 0 4 では data 指示が L2 かどうかを調べる。L2 ならばステップ 1 1 0 6 へ進み、ターゲットキャッ

シュをL1とする。そうでない場合、即ちdata指示がMEMまたは指定なしの場合は、ステップ1105へ進み、ターゲットキャッシュをL2とする。

【0040】

ステップ1107では、タイリングの適用条件を調べ、適用対象かどうかを決定する。適用条件としては、多重ループであるか、依存条件を満たすか、タイリング効果があるかなどを調べる。適用条件の判定方法の詳細については、従来の技術の項で示した非特許文献2に記載されている。適用対象でなければこのループの処理を終了してステップ1101へ戻り、適用対象ならばステップ1108へ進む。

【0041】

ステップ1108では、ステップ1105および1106で決定したターゲットキャッシュのキャッシュサイズ、連想度などに基づいてタイルサイズを決定する。タイルサイズを決定する方法については、従来の技術の項で示した非特許文献2に記載されている。

【0042】

ステップ1109では、ステップ1108で決定したタイルサイズに従い、タイリング変換処理を行う。タイリング変換処理についても、従来の技術の項で示した非特許文献2に記載されている。

【0043】

本実施例では、data指示に応じてターゲットキャッシュを決定することにより、より適切なタイルサイズによる変換を適用する効果がある。また、L1キャッシュにデータがある場合には不要なタイリング変換を適用しない効果がある。

【0044】

次に、本発明によるメモリ最適化502において、ループ交換およびループ展開を行う例について説明する。処理手順を図12に示す。

【0045】

本実施例ではデータがL1キャッシュまたはL2キャッシュにある場合にはループ交換やループ展開の効果がないと考え、適用を中止し、データが主記憶にあると予想される場合のみに最適化を適用する。

【 0 0 4 6 】

ステップ 1 2 0 1 では、未処理のループがあるかどうかを判定し、あればステップ 1 2 0 2 で未処理ループを 1 つ取り出す。ステップ 1 2 0 3 では取り出したループの data 指示をループ表 2 1 0 により調べ、data 指示が L1 または L2 ならば、ステップ 1 2 0 1 へ戻る。これは、このループのデータが既にキャッシュ上にあるため、最適化の効果がないと判断したためである。ステップ 1 2 0 4 ではループ交換の適用条件を調べ、適用対象ならばステップ 1 2 0 5 でループ交換を適用する。ステップ 1 2 0 6 ではループ展開の適用条件を調べ、適用対象ならばステップ 1 2 0 7 でループ展開を適用する。

【 0 0 4 7 】

本実施例ではループ内のデータが主にキャッシュ上に存在する場合には、不要なループ交換やループ展開を適用しない効果がある。

【 0 0 4 8 】

メモリ最適化 5 0 2 において、以上説明したうちの全てまたはいくつかを組み合わせた最適化処理を行ってもよく、いずれか一つの最適化処理を行ってもよい。

【 0 0 4 9 】**【発明の効果】**

本発明によれば、プログラム中のデータアクセスが主としてどのメモリ階層をアクセスするかを指定する手段を設けることにより、データがキャッシュ上に存在するときに効果のない最適化適用の抑止、指定されたメモリ階層に応じたメモリレイテンシを用いた命令スケジューリングやプリフェッチ最適化の適用、または／および指定されたメモリ階層に応じたターゲットキャッシュのキャッシュサイズ等のパラメタを用いたタイリングの適用が可能となり、プログラム実行の高速化が図れる。

【図面の簡単な説明】

【図 1】 本発明の概要を示す図である。

【図 2】 改良された最適化方法を実施するコンパイラが稼動する計算機システムの構成図。

【図 3】 プログラム実行時に主としてアクセスするメモリ階層のコンパイラオプションによる指示の例を示す図。

【図 4】 プログラム実行時に主としてアクセスするメモリ階層の指示文を含むソースプログラムの例を示す図。

【図 5】 コンパイラの処理手順を示す図。

【図 6】 ループ表の例を示す図。

【図 7】 命令スケジューリングの処理手順を示す図。

【図 8】 DAGの一例を示す図。

【図 9】 DAGエッジのレイテンシ設定の処理手順を示す図。

【図 1 0】 プリフェッチの処理手順を示す図。

【図 1 1】 タイリングの処理手順を示す図。

【図 1 2】 ループ交換およびループ展開の処理手順を示す図。

【図 1 3】 データ指示フィールド設定の処理手順を示す図。

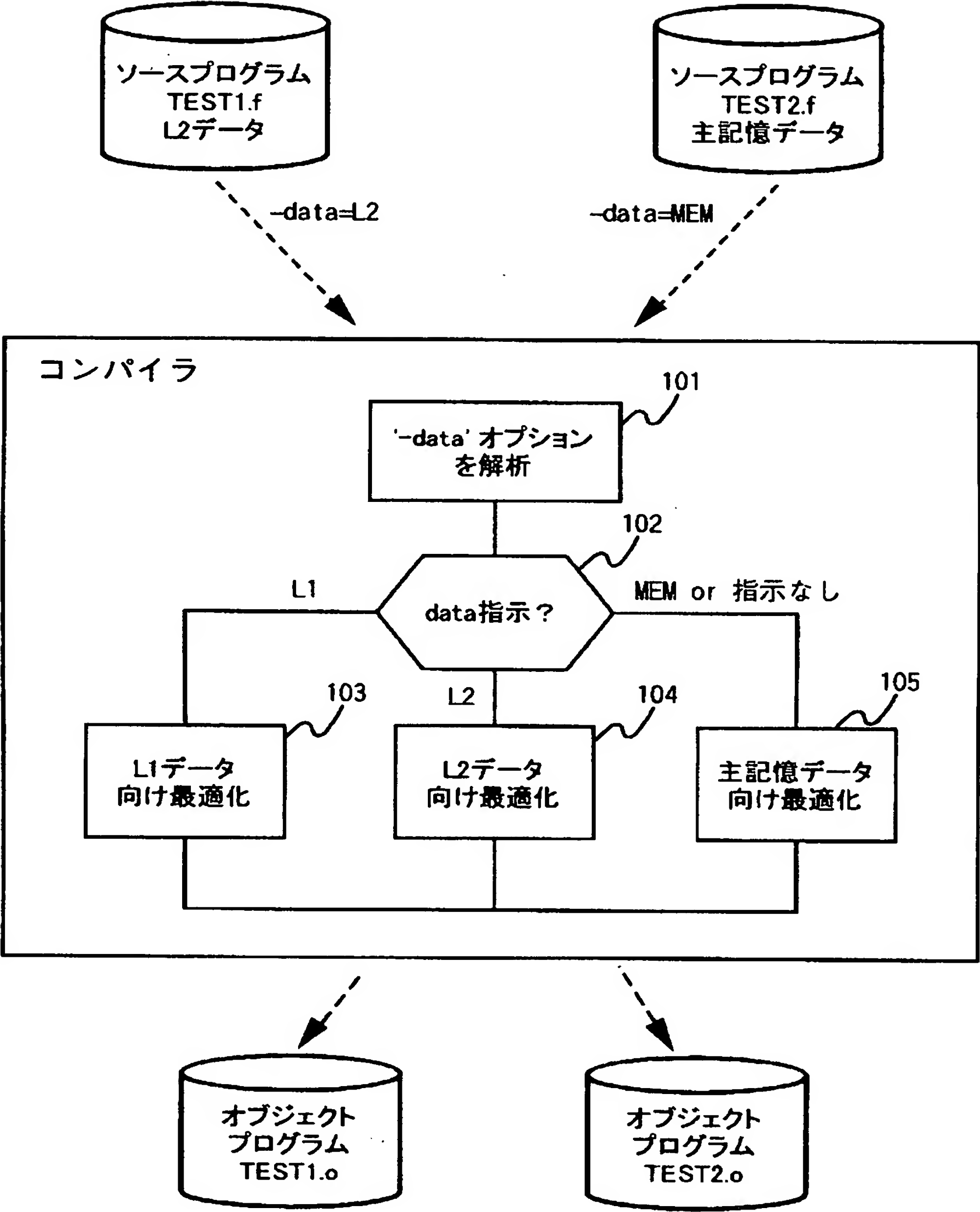
【符号の説明】

2 0 8 . . . コンパイラ、5 0 2 . . . メモリ最適化処理、1 0 3 . . . L1データ向け最適化、1 0 4 . . . L2データ向け最適化、1 0 5 . . . 主記憶データ向け最適化、3 0 1 . . . メモリ階層指示オプション。

【書類名】 図面

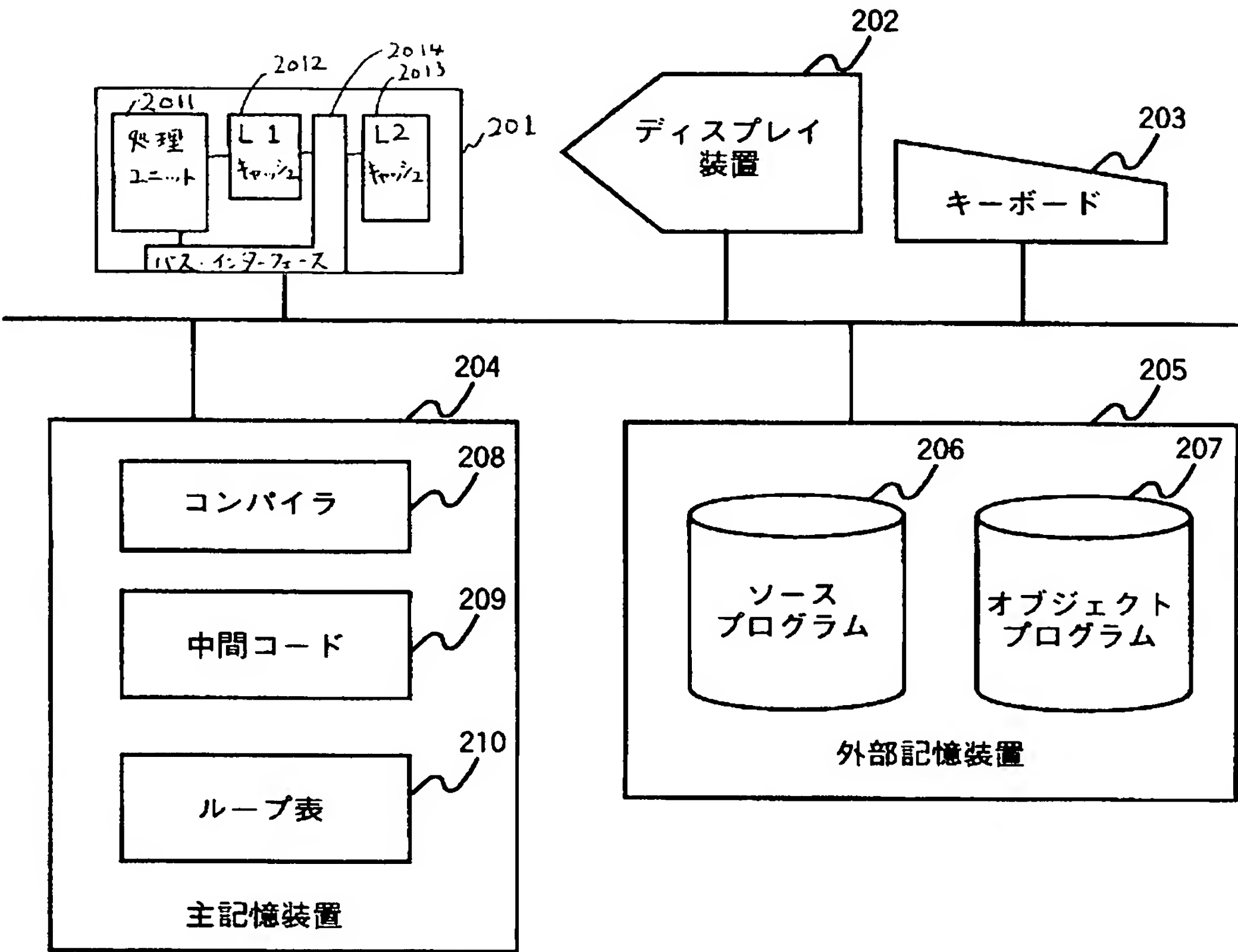
【図 1】

本発明の概要（図 1）



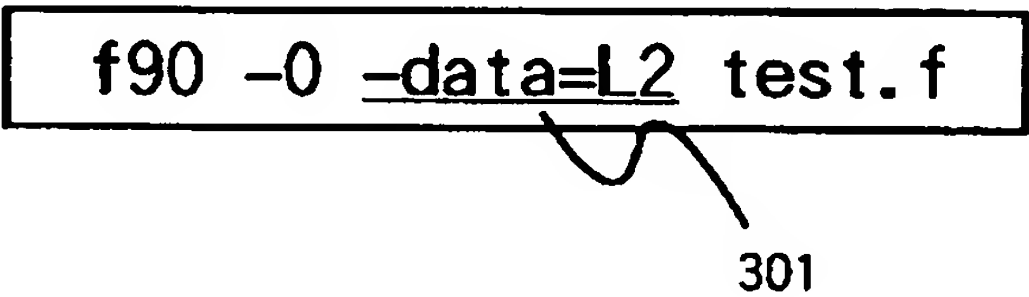
【図 2】

コンパイラが稼動するシステムの構成図（図 2）



【図 3】

コンパイラオプションによる指示の例（図 3）



【図 4】

ユーザー指示付きソースプログラムの例 (図 4)

```
subroutine test1(A,B,C,D,E,F,N)
  real*8 A(N),B(N),C(N),D(N),E(N),F(N)
```

```
  *option data(L1)
```

```
    do 10 I=1,N
      A(I)=B(I)
10 continue
```

401

```
  *option data(MEM)
```

```
    do 20 I=1,N
      C(I)=D(I)
20 continue
```

402

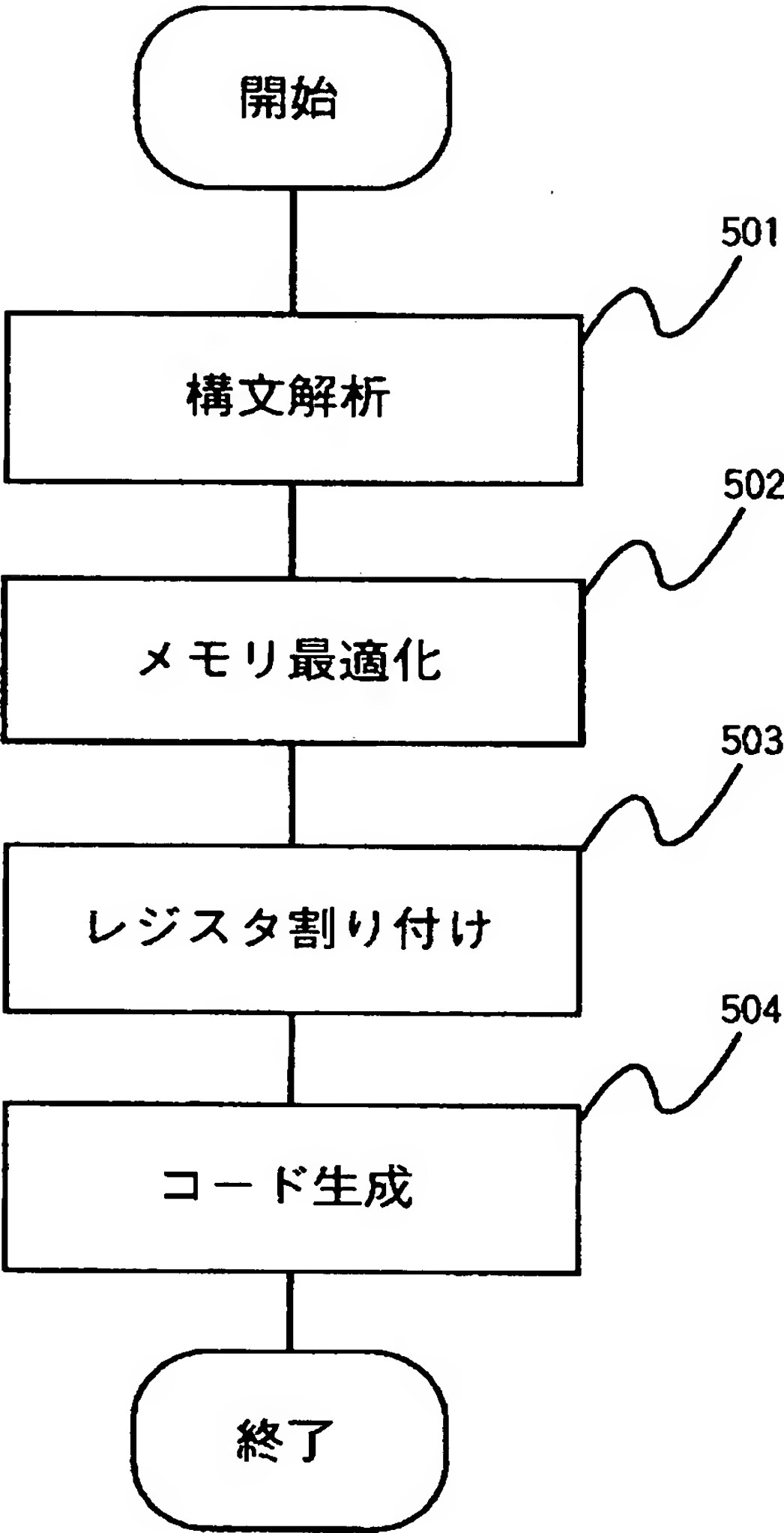
```
  *option data(L2)
```

```
    do 30 I=1,N
      E(I)=F(I)
30 continue
end
```

403

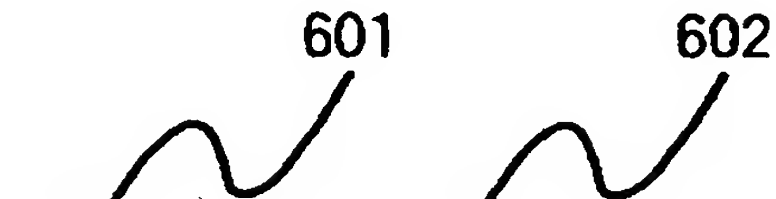
【図 5】

コンパイラの処理手順（図 5）



【図 6】

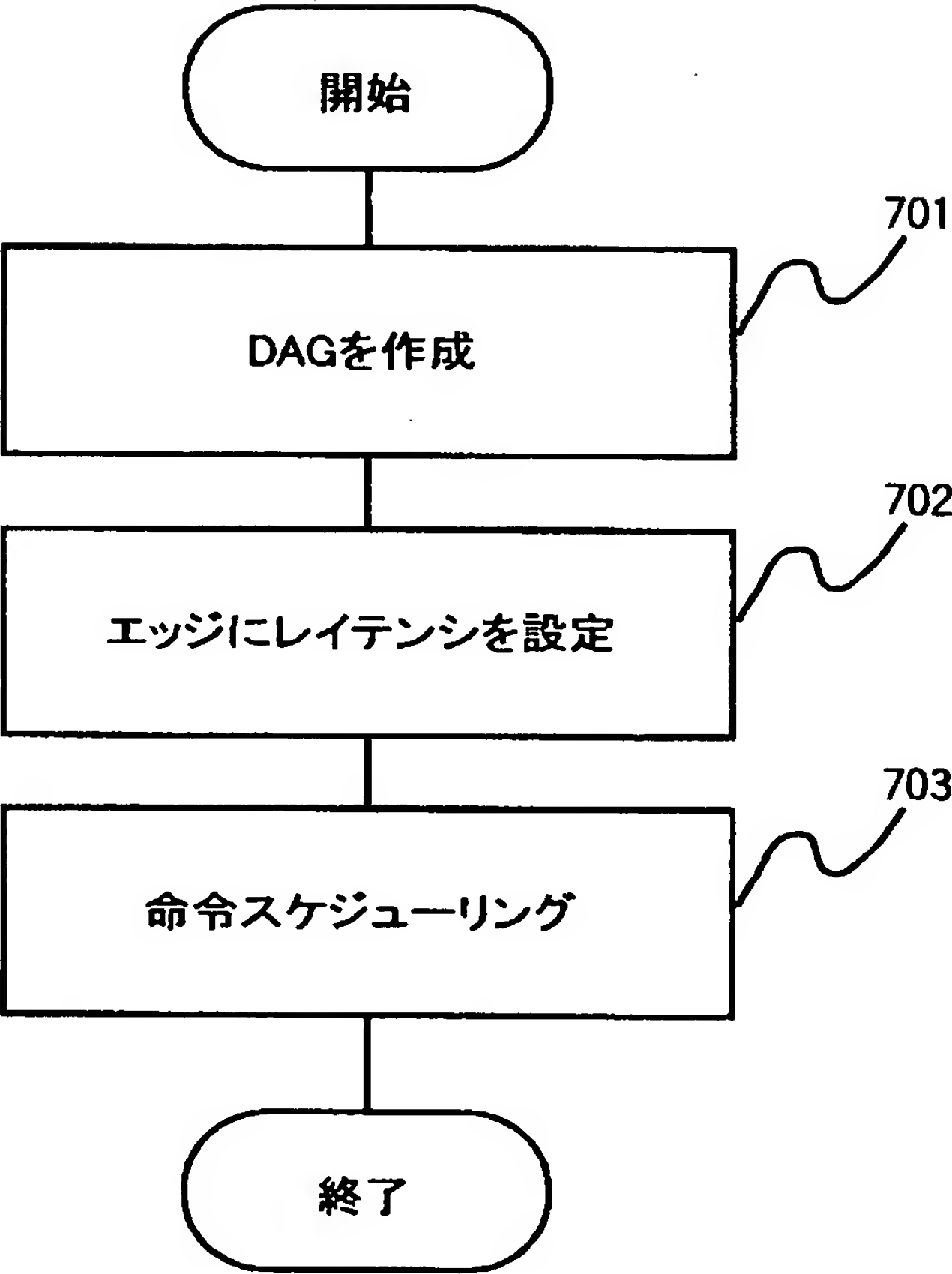
ループ表の例（図 6）



ループ番号	data指示
1	L2
2	MEM
3	指定なし
4	L1

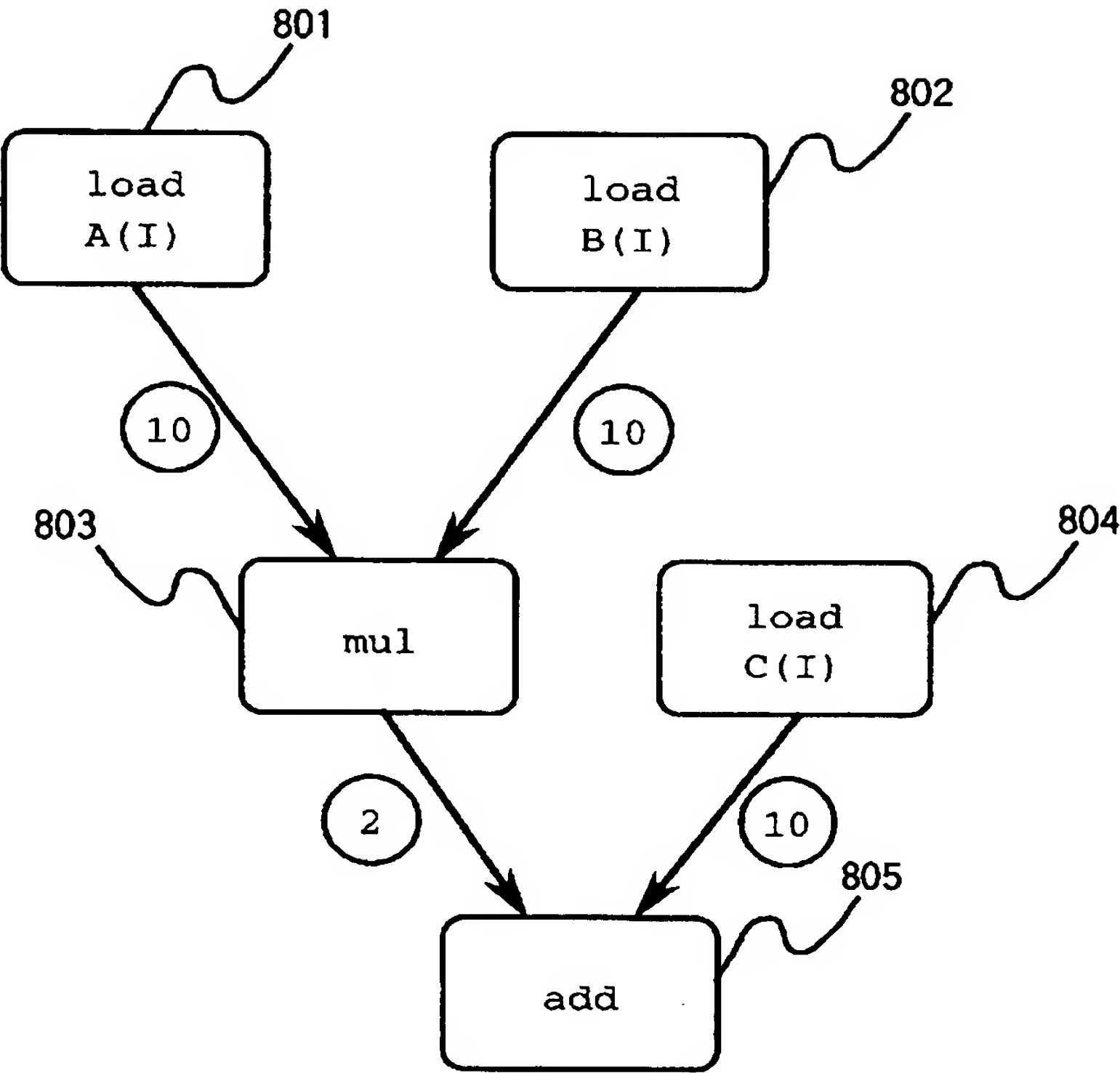
【図 7】

命令スケジューリングの処理手順（図 7）



【図 8】

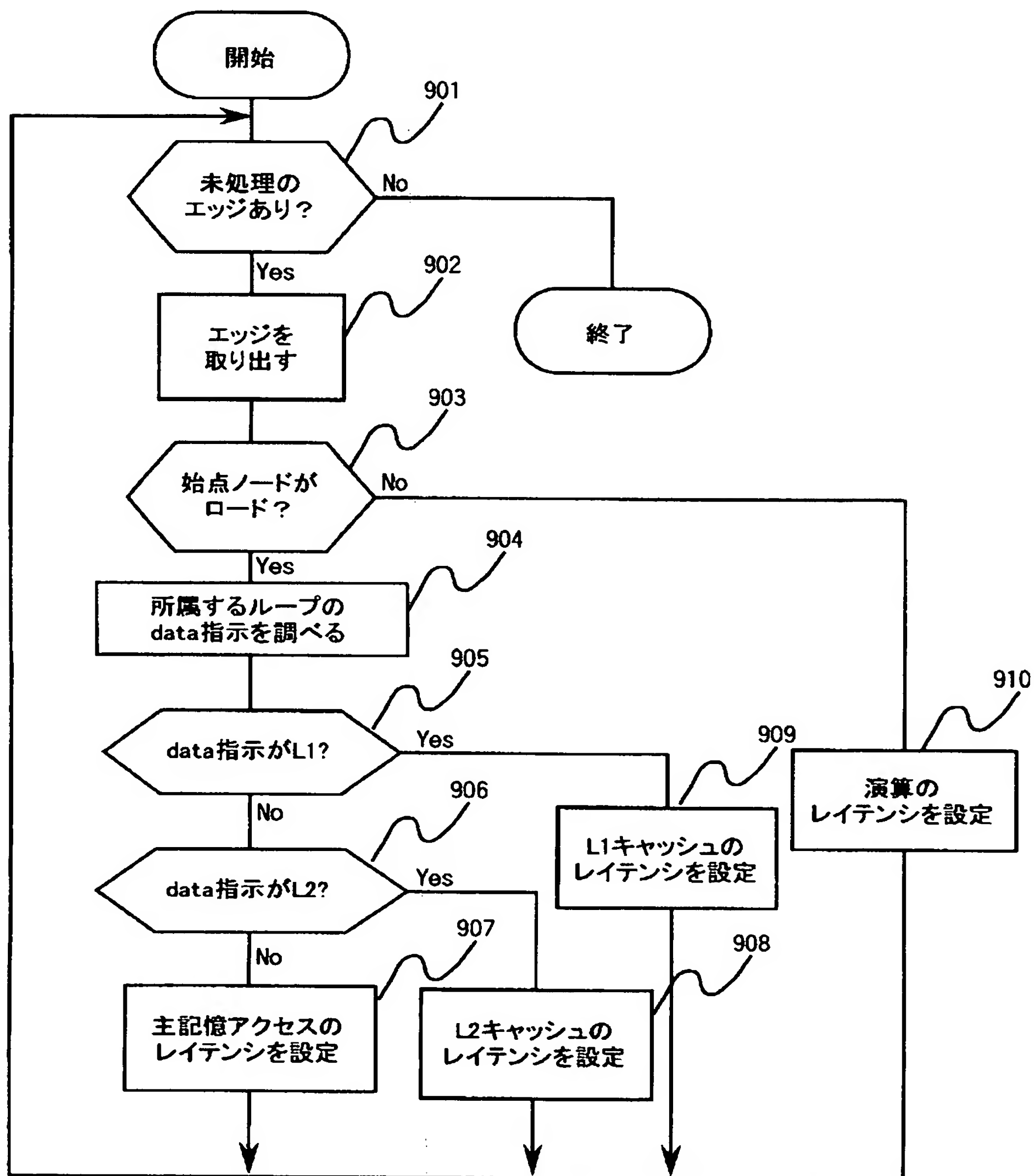
DAG の例（図 8）



○：レイテンシ

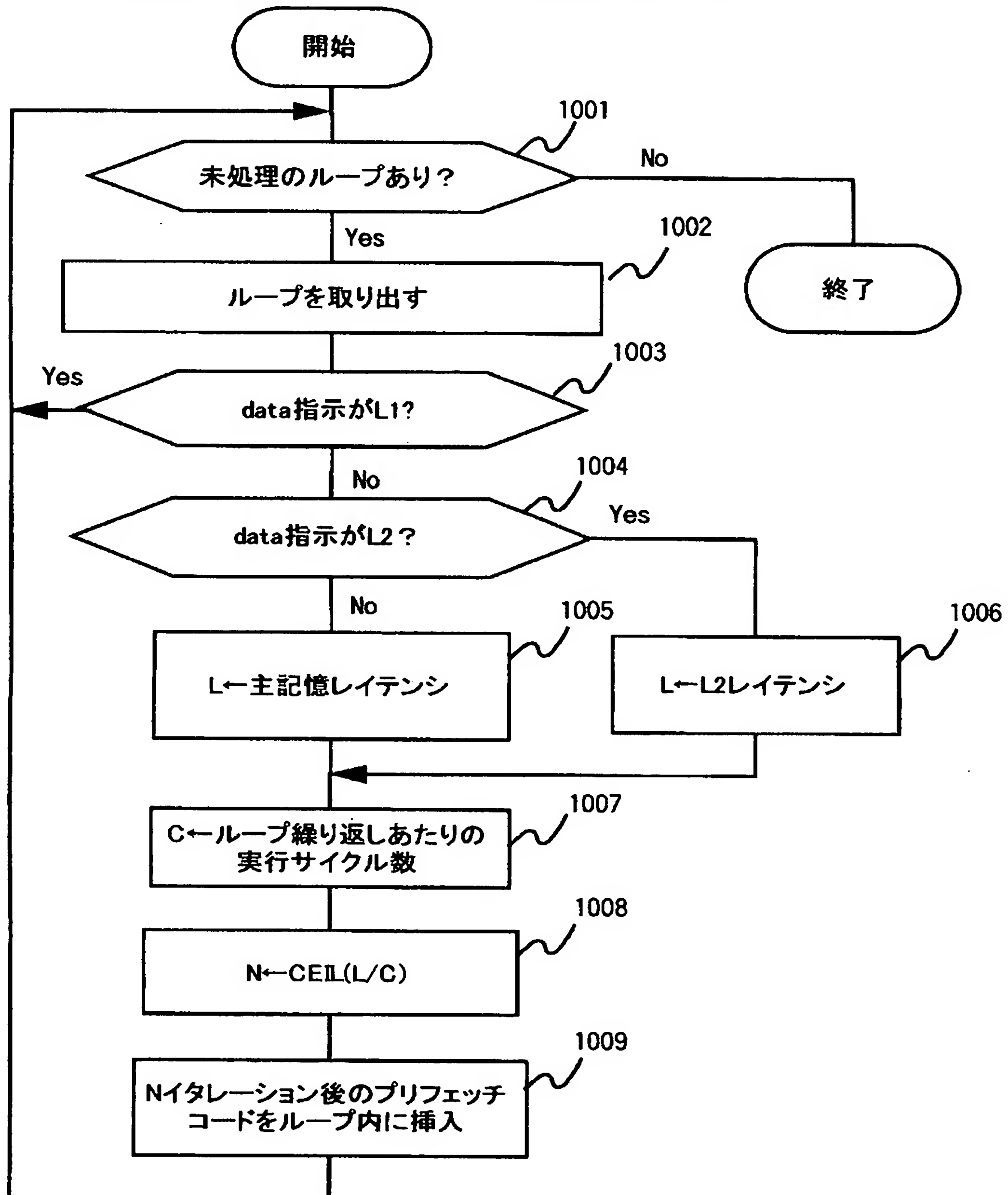
【図 9】

DAGエッジのレイテンシ設定の処理手順（図 9）



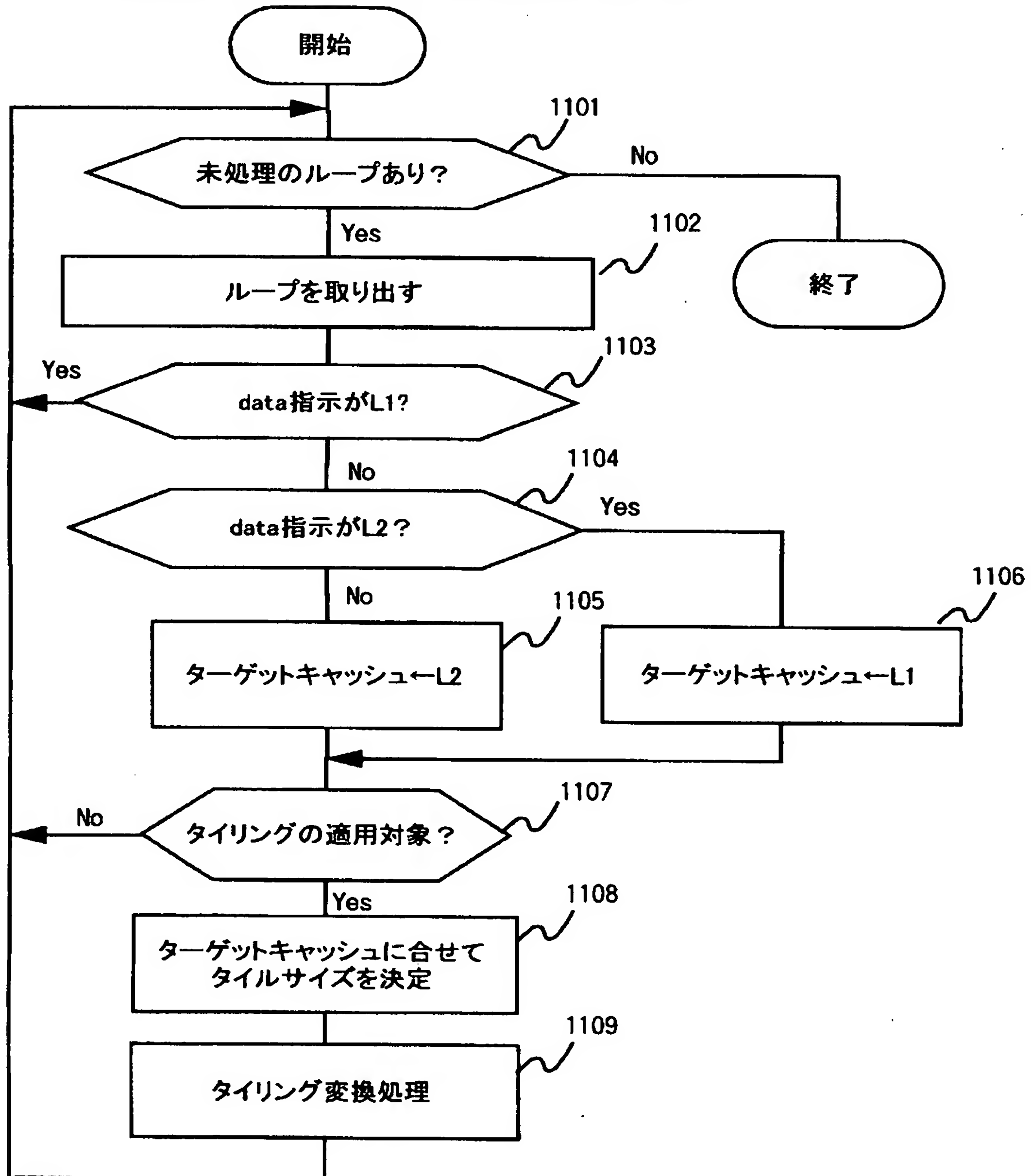
【図 10】

本発明によるプリフェッチの処理手順（図 10）



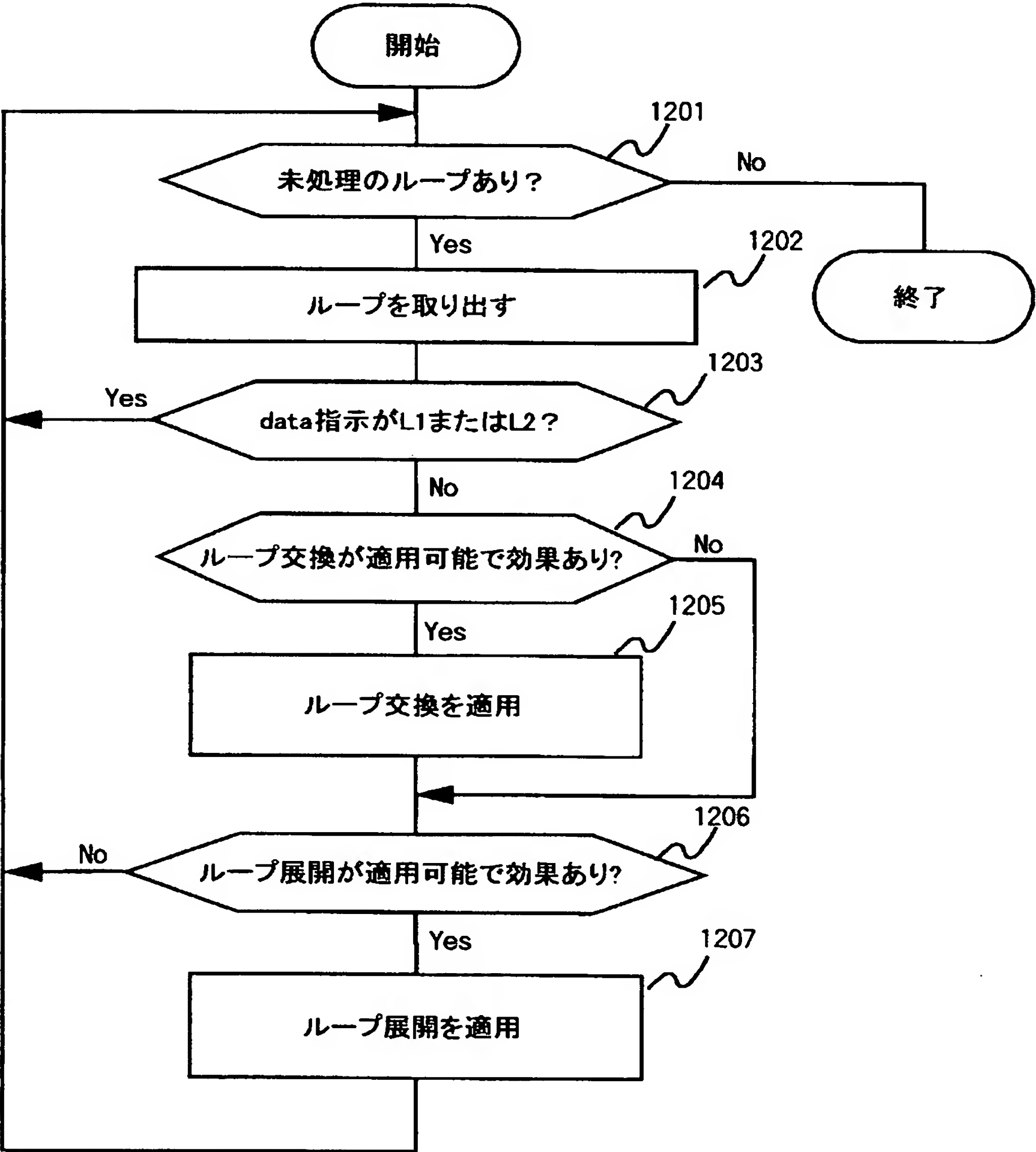
【図 11】

本発明によるタイリングの処理手順（図 11）



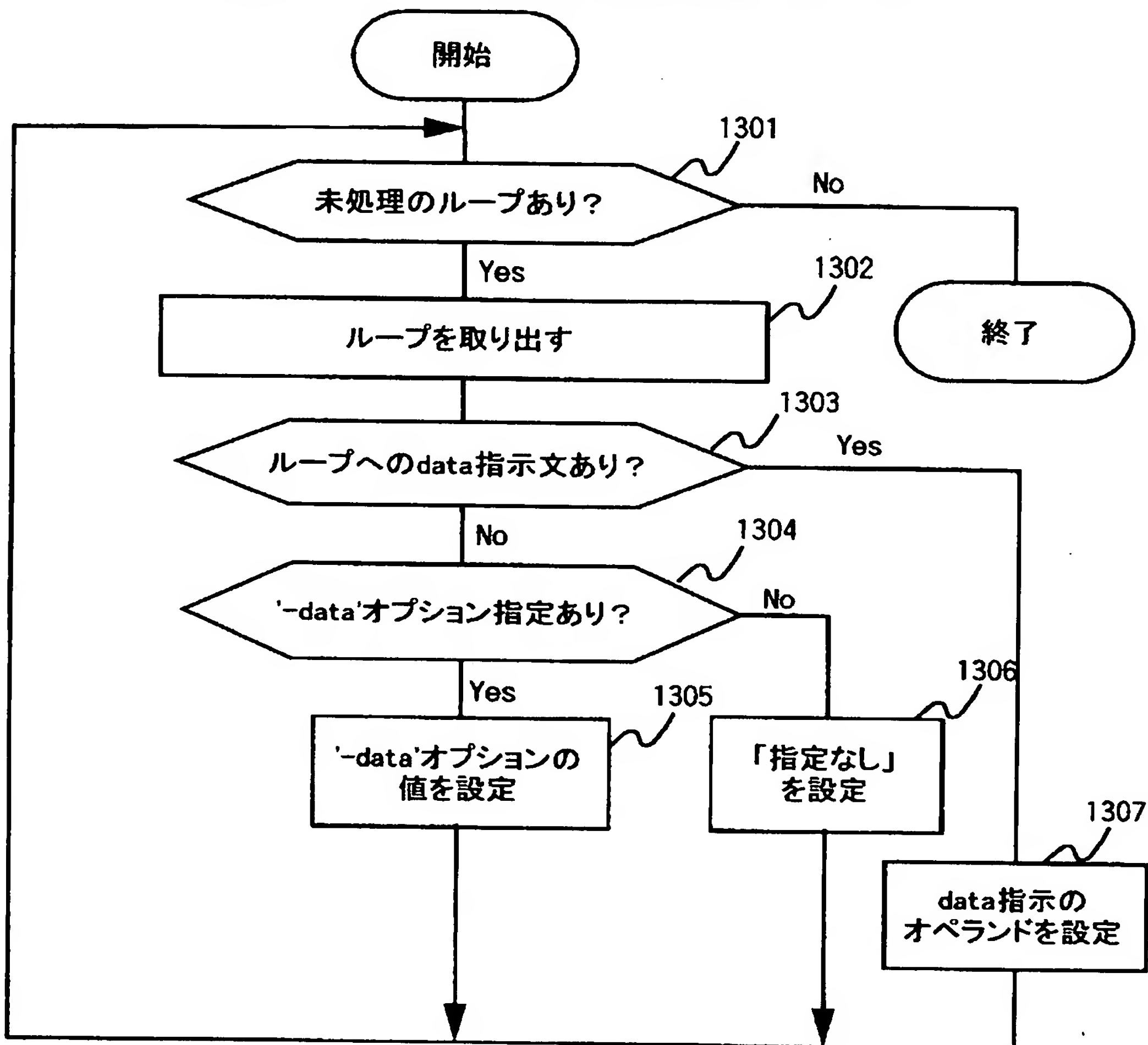
【図 1 2】

本発明によるループ交換・ループ展開の処理手順（図 1 2）



【図 13】

データ指示フィールド設定の処理手順（図 13）



【書類名】 要約書

【要約】

【課題】

プログラムが実行時に主としてアクセスするメモリ階層に応じて、異なる最適化方法を適用する方法を提供する。

【解決手段】

コンパイルオプション指示（3 0 1）またはプログラム中の指示文（4 0 1）を用いて、プログラムが主としてアクセスするメモリ階層をユーザが指定する。コンパイラ（2 0 8）ではメモリ階層指示を解析し（1 0 2）、メモリ階層に応じた最適化（1 0 3、1 0 4、1 0 5）を実施する。

【選択図】 図 1

認定・付加情報

特許出願の番号	特願 2 0 0 3 - 0 9 7 5 7 3
受付番号	5 0 3 0 0 5 3 9 0 1 5
書類名	特許願
担当官	第七担当上席 0 0 9 6
作成日	平成 1 5 年 4 月 2 日

< 認定情報・付加情報 >
【提出日】 平成15年 4月 1日

特願 2 0 0 3 - 0 9 7 5 7 3

出 願 人 履 歴 情 報

識別番号

[0 0 0 0 0 5 1 0 8]

1 . 変更年月日

1 9 9 0 年 8 月 3 1 日

[変更理由]

新規登録

住 所

東京都千代田区神田駿河台 4 丁目 6 番地

氏 名

株式会社日立製作所